# **Effect Analysis of Factors in Generalized Linear Models**

Nobuoki Eshima

Department of Medical Information Analysis, Faculty of Medicine, Oita Medical University, Oita 879-5593, Japan e-mail: eshima@oita-med.ac.jp

## Abstract

Effect analysis of factors on responses or measurements in generalized linear models is important in many application studies, e.g. medicine, sociology, psychology etc. The effect analysis is usually made in testing the main effects and interactions for levels of factors. It is significant to consider summary measures of the effects of factors, i.e. contributions of factors on explaining a response variable. In the present paper, we discuss the partial, association, and total effects of factors in generalized linear models with canonical links. The summary effect measures are provided through a discussion of log odds ratio, and it is shown that the effect measures are related to the Kullback-Leibler information. The present approach is applied to some important generalized linear models.

Key Words: Association effect; Effect analysis; Generalized linear model; Log odds ratio; Partial effect; Total effect.

## 1. Introduction

Analyses of the effects of factors on responses or measurements have been made in many scientific fields, e.g. medicine, sociology, psychology etc. Generalized linear models (Nelder and Wedderburn, 1972; McCullagh and Nelder, 1989) are popularly used for describing regression relationships between response variables and factors in both experimental and observational studies. The generalized linear models (GLMs) are an exponential family of distributions that include various important regression models, e.g. ordinary linear regression model (Draper and Smith, 1966), logistic regression (logit) model (Cox and Snell, 1989; Hosmer and Lemeshow, 1989), Poisson regression model (Vonesh, 1990; Christiansen and Moris, 1997), Loglinear model (Haberman, 1979) etc. Generalized linear models are also applied to statistical quality control (Hamada and Nelder, 1997), because there are many cases of nonnormal responses in industrial experiments (Lewis, et al. 2001a). The usual least square approach with response transformation and the GLM approach are compared, and Lewis et al. (2001a, b) investigated advantages of the latter approach.

The usual effect analysis with GLMs is to statistically test the main effects and the interactions of levels of factors, and to interpret the results. After the analysis, it is important to summarize the effects of the factors, that is, comparison of contributions of factors on the variation of response *Y*. In this paper, the partial, association and total effects of factors are discussed for generalized linear models. The effect measures for the partial, association and total effects are derived from a discussion of log odds ratio, and properties of the effect measures are considered. Examples of

GLMs are also provided.

## 2. Basic Idea of Measuring the Effects of Factors in GLMs with Canonical Links

Let  $X = (X_1, X_2, ..., X_K)^T$  be a  $K \times 1$  factor vector; let Y be a response variable; and let  $f(y|\mathbf{x})$  be the conditional probability or density function of Y given  $\mathbf{X} = \mathbf{x} = (x_1, x_2, ..., x_K)^T$ . Then, a GLM with the canonical link is described as follows:

$$f(y|\mathbf{x}) = \exp\{(y - b(y))/a(y) + c(y, y)\},$$
(2.1)

where and are parameters and for  $= \begin{pmatrix} 1 & 2 \\ 1 & 2 \end{pmatrix}^{T}$  $= {}^{T} \boldsymbol{x}.$  (2.2)

Without loss of generality, we set a() > 0. The effect of factor  $X_i = x_i$  may be defined by i from an analogy to the ordinary linear regression model, because the predictor (2.2) increases by i for unit change of the factor. However, it is a question how this effect is interpreted except the usual regression model. Let  $OR(y,y^*|x,x^*)$  be the odds ratio with respect to  $x, x^*, y$ , and  $y^*$ . Then, we have  $\log OR(y,y^*|x,x^*) = (y - y^*)^{-T}(x - x^*)/a(-)$ . (2.3)

The first interpretation of this quantity is that (2.3) is the inner product of the response  $y - y^*$  and the predictor  ${}^{T}(x - x^*)$  with respect to a(). The second interpretation is as follows. The conditional log odds of Y = y over  $Y = y^*$  given X = x is

 $\log\{f(y|\mathbf{x})/f(y^*|\mathbf{x})\} = \{-\log f(y^*|\mathbf{x})\} - \{-\log f(y|\mathbf{x})\}$ 

= the amount of the conditional uncertainty of  $Y = y^*$  minus that of Y = y given X = x.

From this, the above quantity is a decrease of uncertainty of y over  $y^*$ . Hence, the log odds ratio (2.3) can be interpreted as the change of uncertainty of y over  $y^*$  for the change of factors from  $x^*$  to x. When we substitute the baselines  $x^*$  and  $y^*$  in (2.3) for the expectations  $\mu_X$  and  $\mu_Y$ , respectively, we get

$$\log OR(y, \boldsymbol{\mu}_{\boldsymbol{X}} | \boldsymbol{x}, \boldsymbol{\mu}_{\boldsymbol{X}}) = (y - \boldsymbol{\mu}_{\boldsymbol{Y}})^{-T} (\boldsymbol{x} - \boldsymbol{\mu}_{\boldsymbol{X}})/a(-).$$
(2.4)  
dds ratio can be viewed as the total effect of  $\boldsymbol{X} - \boldsymbol{x}$  on  $\boldsymbol{X} - \boldsymbol{y}$ . In order to summarize the above

This odds ratio can be viewed as the total effect of X = x on Y = y. In order to summarize the above quantity, we define the total effect of X on Y by

$$e_{T}(X \quad Y) = E\{(Y - \mu_{Y}) \quad {}^{T}(X - \mu_{X})/a()\}.$$
(2.5)

This quantity is the covariance of the response *Y* and the predictor  ${}^{T}X$  with respect to a(). For example, the following linear regression model is considered:

$$Y = \mathbf{\mu} + \mathbf{x} + e, \qquad (2.6)$$

where error term *e* is distributed according to normal distribution  $N(0, ^2)$ . In this case, the dispersion parameter is  $a() = ^2$ . From this model, we have

 $e_{T}(X \ Y) = E\{(X - \mu_{X})^{T} \ ^{T}(X - \mu_{X})\}/^{2} = (Var(Y) - ^{2})/^{2}.$ 

This is the ratio of the explained variation of Y by factor vector X to the variation of the error term. Since  $a(\ )$  in (2.5) is a dispersion parameter of the GLM (2.1), the parameter is referred to as a generalized dispersion of the GLM (2.1), i.e. the generalized dispersion of the error of the prediction. The numerator in (2.5)

$$\mathrm{E}\{(Y - \boldsymbol{\mu}_{Y}) \quad ^{\mathrm{T}}(X - \boldsymbol{\mu}_{X})\}$$

can be viewed as a generalized explained dispersion of response Y by the predictor vector X. Hence,

the total effect (2.5) can be interpreted as the ratio of a generalized explained variation of response Y by factor X to a generalized error variation of the GLM (2.1).

Second, we consider the partial effects of  $X_i$  on Y. Let  $X^{(i)} = (X_1, X_2, ..., X_{i-1}, X_{i+1}, ..., X_K)^T$ ; let  $\mu$ 

 $_{i}(\mathbf{x}^{(i)})$  be the conditional expectation of  $X_{i}$  given  $\mathbf{X}^{(i)} = \mathbf{x}^{(i)}$ ; let  $\mathbf{\mu}_{X}(\mathbf{x}^{(i)}) = (x_{1}, x_{2}, ..., x_{i-1}, \mathbf{\mu}_{i})$  $_{i}(\mathbf{x}^{(i)}), x_{i+1}, ..., x_{K})^{\mathrm{T}}$ , i.e. the conditional expectation of X given  $\mathbf{X}^{(i)} = \mathbf{x}^{(i)}$ ; and let  $\mathbf{\mu}_{Y}(\mathbf{x}^{(i)})$  be the conditional expectation of Y given  $\mathbf{X}^{(i)} = \mathbf{x}^{(i)}$ . Then, we get

log OR(y,  $\mu_{Y}(\mathbf{x}^{(i)})|\mathbf{x}, \mu_{X}(\mathbf{x}^{(i)})) = (y - \mu_{Y}(\mathbf{x}^{(i)})) _{i}(x_{i} - \mu_{i}(\mathbf{x}^{(i)}))/a()$ . (2.7) The quantity (2.7) can be viewed as the direct effect of level  $X_{i} = x_{i}$  on Y = y. By taking the conditional expectation with respect to  $X_{i}$  and Y given  $\mathbf{X}^{(i)} = \mathbf{x}^{(i)}$ , we have

 $E\{\log OR(y, \boldsymbol{\mu}_{\boldsymbol{Y}}(\boldsymbol{x}^{(i)})|\boldsymbol{x}, \boldsymbol{\mu}_{\boldsymbol{X}}(\boldsymbol{x}^{(i)})|\boldsymbol{X}^{(i)} = \boldsymbol{x}^{(i)})\} = iCov(\boldsymbol{Y}, X_i | \boldsymbol{X}^{(i)} = \boldsymbol{x}^{(i)})/a(), \quad (2.8)$ where  $Cov(\boldsymbol{Y}, X_i | \boldsymbol{X}^{(i)} = \boldsymbol{x}^{(i)})$  are the conditional covariance of  $\boldsymbol{Y}$  and  $X_i$  given  $\boldsymbol{X}^{(i)} = \boldsymbol{x}^{(i)}$ . The partial effect of  $X_i$  on  $\boldsymbol{Y}$  is defined by the expectation of the above quantity with respect to  $\boldsymbol{X}^{(i)}$ :

$$p_{\rm P}(X_i \quad Y) = {}_i {\rm Cov}(Y, X_i | X^{(i)}) / a(\ ), \qquad (2.9)$$

where  $\text{Cov}(Y,X_i|X^{(i)})$  is the expectation of  $\text{Cov}(Y,X_i|X^{(i)}=x^{(i)})$  with respect to  $X^{(i)}$ , i.e. for the marginal density or probability function of  $X^{(i)}$ ,  $g(x^{(i)})$ ,

$$\operatorname{Cov}(Y,X_i|X^{(i)}) = \operatorname{Cov}(Y,X_i|X^{(i)} = x^{(i)})g(x^{(i)})dx^{(i)}.$$

The partial effect (2.7) can be interpreted as the ratio of the generalized partially explained variation of response *Y* by factor  $X_i$  to a generalized error variation of the GLM (2.1). For example, in an ordinary linear regression model (2.6) we have

$$\mathbf{e}_{\mathbf{P}}(X_i \quad Y) = \frac{2}{i} \operatorname{Var}(X_i \mid \mathbf{X}^{(i)}) / \frac{2}{i}$$

Third, the total effects of factors are discussed. According to the above discussion, the direct effect of factor vector  $X^{(i)}$  is defined by

$$e_{P}(X^{(i)} \quad Y) = \int_{j=i}^{j} Cov(Y, X_{j} | X_{i})/a().$$
 (2.10)

The total effect of factor  $X_i$  is defined by

 $\mathbf{e}_{\mathrm{T}}(X_i \quad Y) = \mathbf{e}_{\mathrm{T}}(X \quad Y) - \mathbf{e}_{\mathrm{P}}(X^{(i)} \quad Y)$ 

$$= {}_{i} \operatorname{Cov}(Y, X_{i}) / a( ) + {}_{j i} \operatorname{Cov}(\mu_{j}(X_{i}), \mu_{Y}(X_{i})) / a( ).$$
(2.11)

The total effect is interpreted as the inner product of  $y - \mu_{y}$  and the predictor

$$(x_i - \mu_i)/a() + j(\mu_j(x_i) - \mu_j)/a().$$

By subtracting the partial effect of  $X_i$  (2.9) from the total effect (2.11), we define the association effect of  $X_i$  by

$$e_{A}(X_{i} \quad Y) = e_{T}(X_{i} \quad Y) - e_{P}(X_{i} \quad Y) = {}_{i}Cov(\mu_{i}(X^{(i)}), \mu_{Y}(X^{(i)}))/a() + {}_{j \quad i} {}_{j}Cov(\mu_{j}(X_{i}), \mu_{Y}(X_{i}))/a()).$$
(2.12)

The above effect is interpreted as the effect that is made by the association between  $X_i$  and  $X^{(i)}$ . If  $X_i$  is independent of  $X^{(i)}$ ,

$$e_{A}(X_{i} \qquad Y) = 0.$$

The total effect is

$$e_{\mathrm{T}}(X_i \qquad Y) = e_{\mathrm{P}}(X_i \qquad Y) + e_{\mathrm{A}}(X_i \qquad Y).$$

By replacing the factor  $X_i$  in the above discussion by the set of factors, the total, direct and indirect effects of the set of factors can be defined. The following quantity is the contribution proportion of the effect of factor  $X_i$  on Y:

$$C(X_i) = e_T(X_i \qquad Y)/e_T(X \quad Y).$$
(2.13)

It is seen that the effects defined above are scale-invariant. If  $X_i$  (i = 1, 2, ..., K) are mutually independent, from (2.5) we have

$$e_{T}(X \ Y) = \sum_{i=1}^{K} e_{P}(X_{i} \ Y) = \sum_{i=1}^{K} e_{T}(X_{i} \ Y).$$
 (2.14)

From the above consideration, we have

$$e_{\mathrm{T}}(\boldsymbol{X} \qquad \boldsymbol{Y}) = e_{\mathrm{T}}(X_i \qquad \boldsymbol{Y}) + e_{\mathrm{P}}(\boldsymbol{X}^{(i)} \qquad \boldsymbol{Y})$$
$$= e_{\mathrm{P}}(X_i \qquad \boldsymbol{Y}) + e_{\mathrm{T}}(\boldsymbol{X}^{(i)} \qquad \boldsymbol{Y}).$$

## **3.** Application to GLMs

In this section, the above idea is applied to normal distributions and the logit model. The details are omitted.

## 4. Numerical Examples

The data concerning the effects of AZT in showing the development of AIDS sysamptoms (Agresti, 1996, p.119) and the data in a study of length of time spent on individual home visit by public nurses (Daniel, 1999, pp. 348-353) are reanalyzed with the present approach. The details are omitted.

#### REFERENCES

Agresti, A. (1990). Categorical Data Analysis, New York: John Wiley & Sons, Inc.

- Agresti, A. (1996). An Introduction to Categorical Data Analysis, New York: John Wiley & Sons, Inc.
- Christiansen, C. L. and Morris, C. N. (1997). Hierarchical Poisson regression modeling, *Journal of the American Statistical Association*, **92**, 618-632.
- Cox, D. R. and Snell, E. J. (1989). Analysis of Binary Data, London: Chapman and Hall.
- Daniel, W. W. (1999). *Biostatistics: A Foundation for Analysis in the Health Sciences*, Seventh Edition, New York; John Wiley & Sons, Inc.
- Drayper, N. R. and Smith, H. (1966). *Applied Regression Analysis*, New York: John Wiley & Sons, Inc.
- Haberman, S. J. (1979). Analysis of Qualitative Data, Vol. 2, New York: Academic Press, Inc.
- Hamada, M. and Nelder, J. A. (1997). Generalized linear models for quality-improvement experiments, *Journal of Quality Technology*, 29, 292-304.
- Hosmer, D. and Lemeshow, S. (1990). *Applied Logistic Regression*, New York; John Wiley & Sons, Inc.
- Lewis, S. L., Montgomery, D. C., and Myers, R. H. (2001). Confidence interval coverage for designed experiments analyzed with generalized linear models, *Journal of Quality Technology* (to appear).
- Lewis, S. L., Montgomery, D. C., and Myers, R. H. (2001). Examples of designed experiments with nonnormal responses, *Journal of Quality Technology* (to appear).
- McCullagh, P. and Nelder, J.A. (1989) *Generalized Linear models*, 2<sup>nd</sup> ed. London: Chapman and Hall.
- Nelder, J. A. and Wedderburn, R.W.M. (1972) Generalized linear model, *Journal of the Royal Statistical Society*, A, 135, 370-384.
- Vonesh, E. F. (1990). Modelling Peritonitis rates and associated risk factors for individuals on continuous ambulatory dialysis, *Statistics in Medicine*, 9, 263-271.